# Convolutional Neural Nets
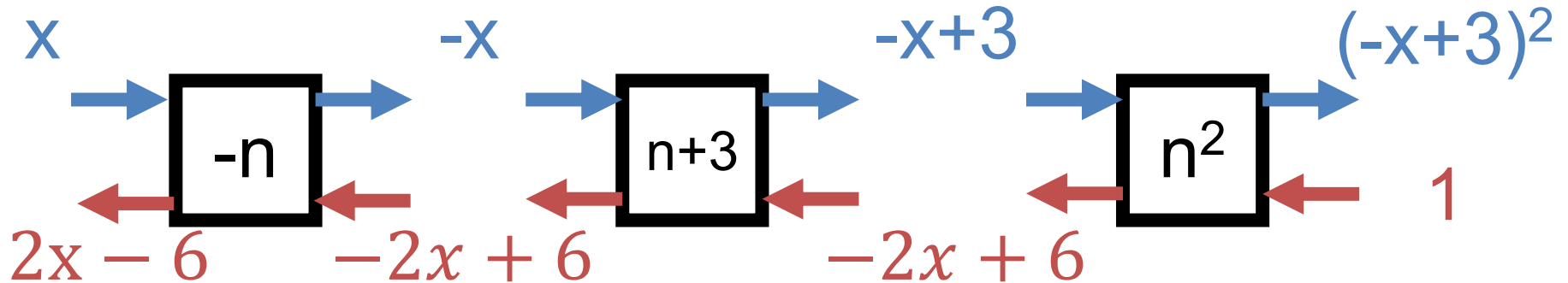
EECS 442 – Jeong Joon Park

Winter 2024, University of Michigan

# Administrivia

- HW3 Due this Wednesday

# Previously – Backpropagation

$$f(x) = (-x + 3)^2$$

x           -x           -x+3        $(-x+3)^2$

| -n | | n+3 | | $n^2$ | 1 |

$2x - 6$     $-2x + 6$     $-2x + 6$

**Forward pass: compute function**
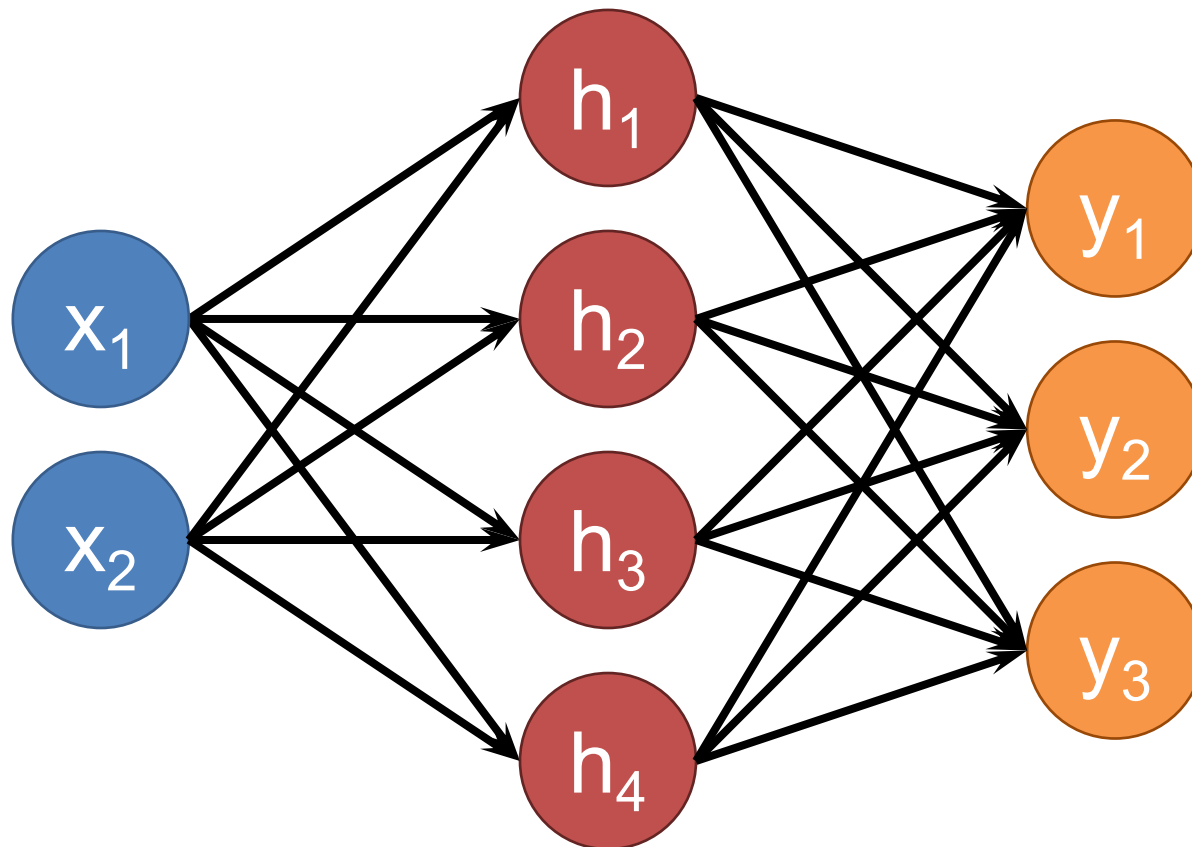
**Backward pass: compute derivative of all parts of the function**

# Setting Up A Neural Net
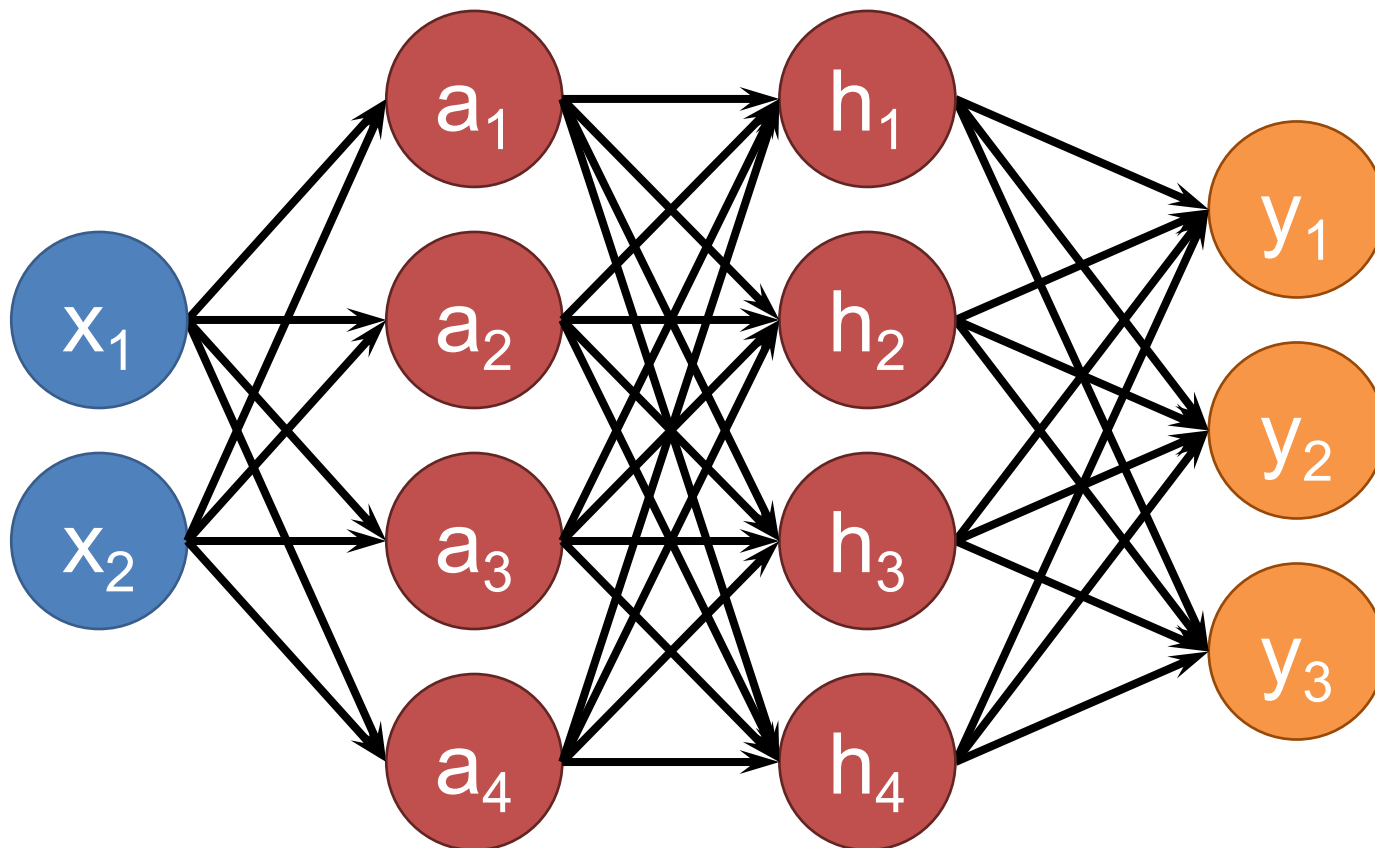
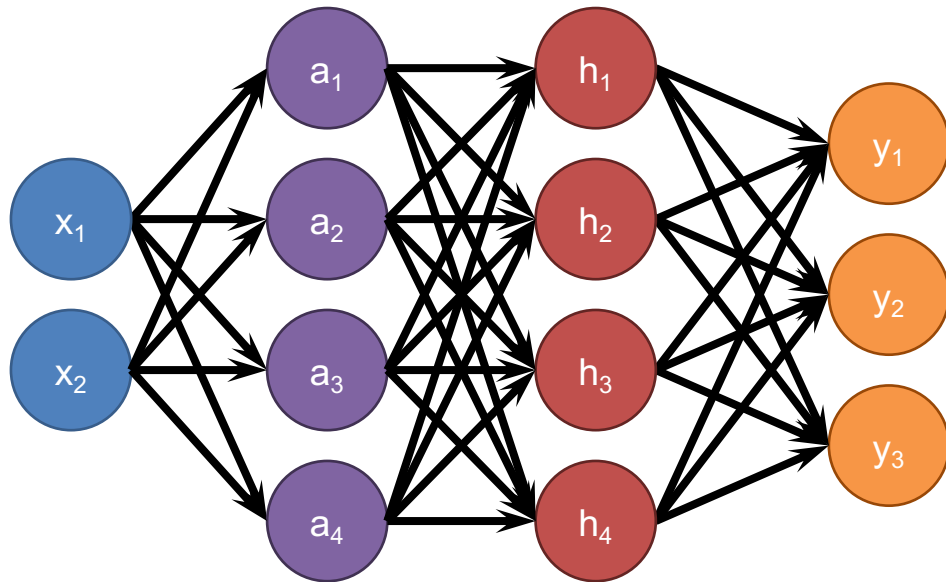Input             Hidden              Output

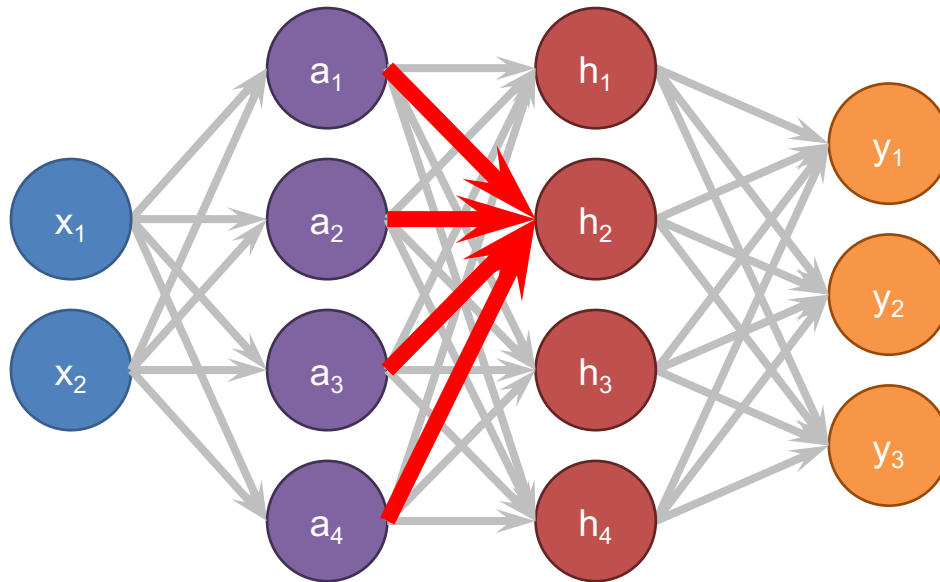# Setting Up A Neural Net

Input  Hidden 1  Hidden 2  Output

# Fully Connected Network



Each neuron connects to each neuron in the previous layer
Fully-Connected Net
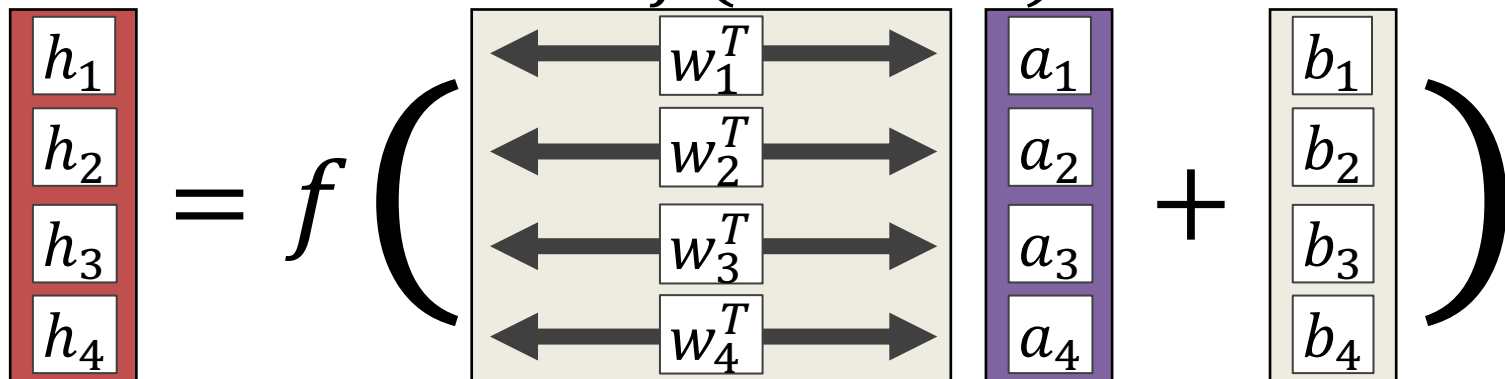
# Fully Connected Network



$a$    All layer *a* values

$w_i, b_i$   Neuron i weights, bias

$f$    Activation function
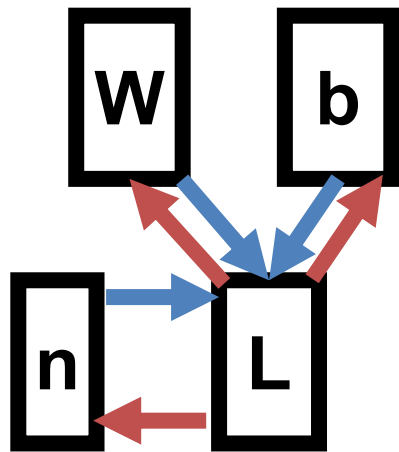
$$h = f(Wa + b)$$

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \end{bmatrix} = f\left( \begin{bmatrix} \longleftrightarrow w_1^T \longleftrightarrow \\ \longleftrightarrow w_2^T \longleftrightarrow \\ \longleftrightarrow w_3^T \longleftrightarrow \\ \longleftrightarrow w_4^T \longleftrightarrow \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \right)$$

# Fully Connected Network

## Define New Block: "Linear Layer"
### *(It's Technically Affine)*

$$L(\boldsymbol{n}) = \boldsymbol{W}\boldsymbol{n} + \boldsymbol{b}$$
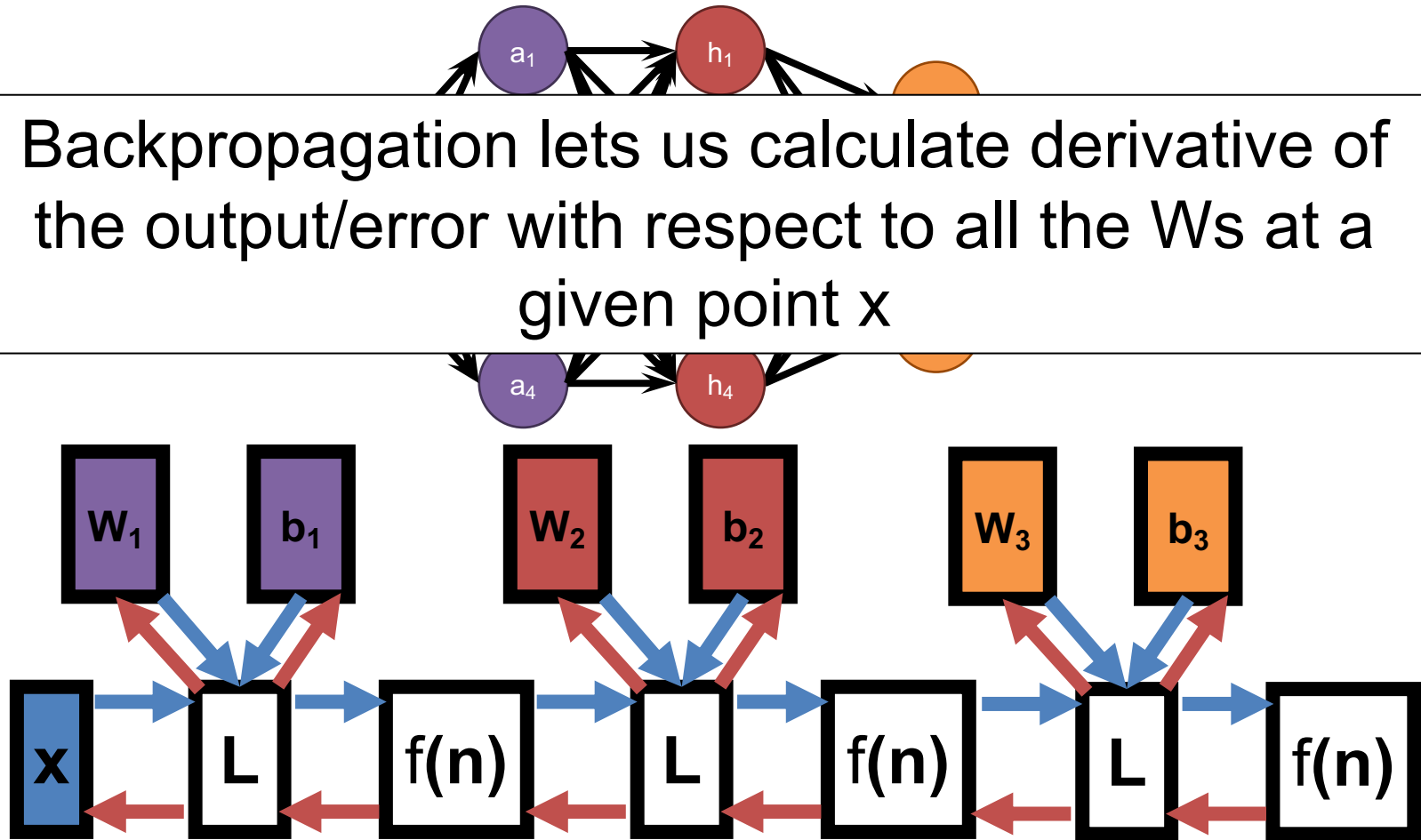
Can get gradient with respect to all the inputs
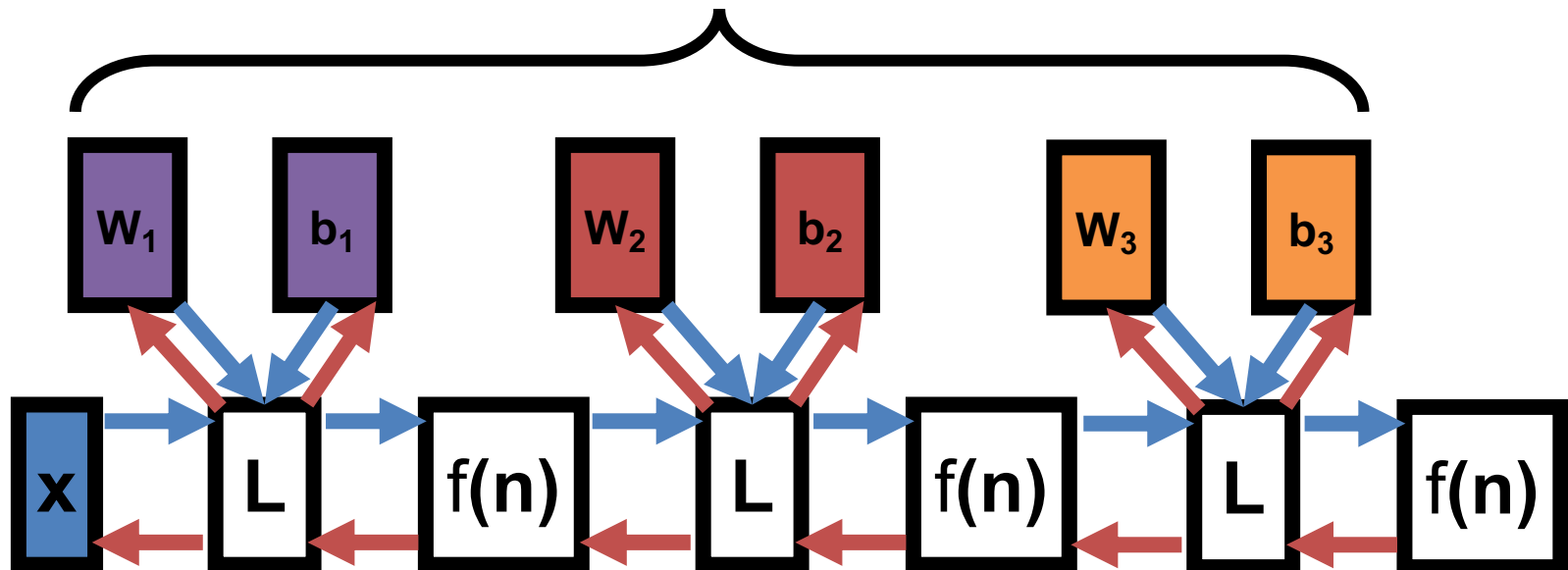
# Fully Connected Network

# Fully Connected Network



Backpropagation lets us calculate derivative of the output/error with respect to all the Ws at a given point x

# Putting It All Together – 1
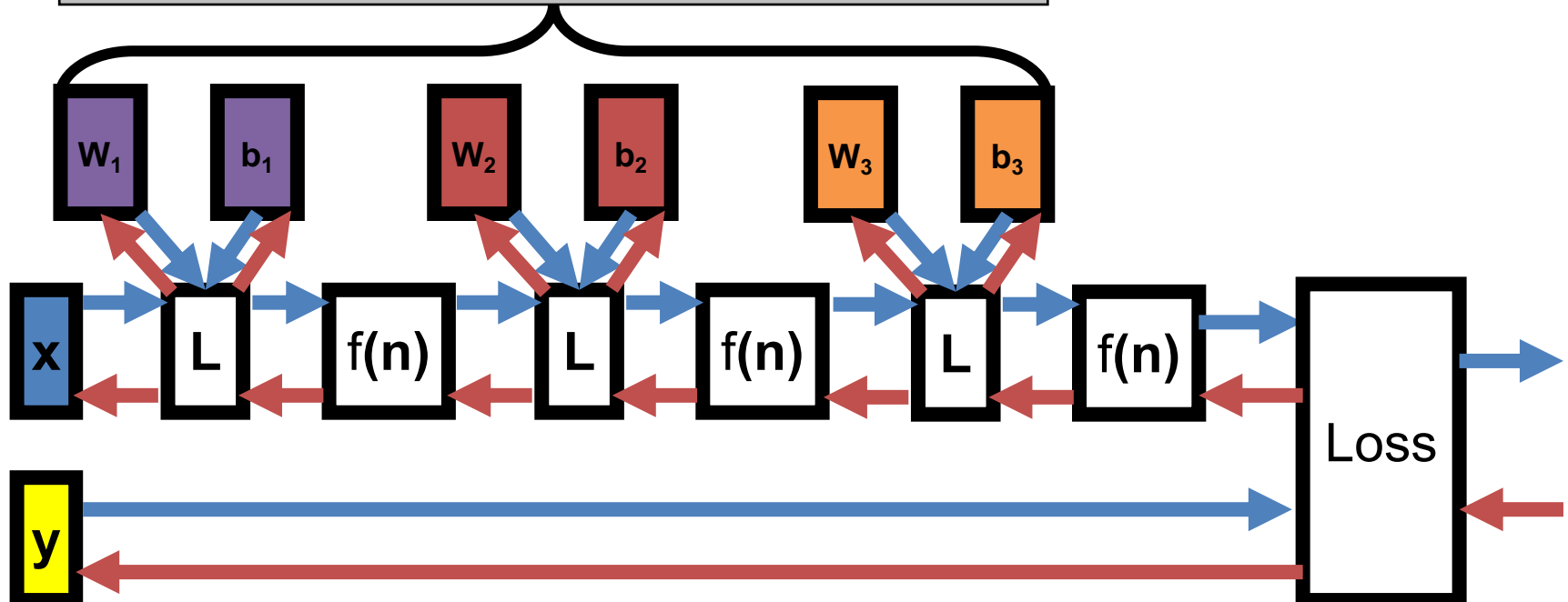
Function: $NN(x; W_i, b_i)$
Parameterized by $W = \{W_i, b_i\}$

# Putting It All Together



Function: Loss(NN(x; $W_i$, $b_i$), y)

Function: NN(x; $W_i$, $b_i$)

$W_1$ $b_1$ $W_2$ $b_2$ $W_3$ $b_3$

x L f(n) L f(n) L f(n) Loss
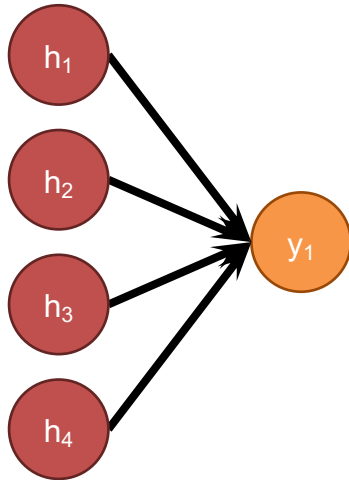
y

# Putting It All Together

```
W = initializeWeights()
for i in range(numIterations):
        #sample a batch
        batch = random.subset(0,#datapoints,K)
        batchX, batchY = dataX[batch], dataY[batch]
        #compute gradient with batch
        gradW = backprop(Loss(NN(batchX,W),batchY))
        #update W with gradient step or use momentum
        W += -stepsize*gradW
return W
```

# What Can We Represent?



$$L(\boldsymbol{n}) = \boldsymbol{W}\boldsymbol{n} + \boldsymbol{b}$$

# What Can We Represent

- Recall: ax+by+z is
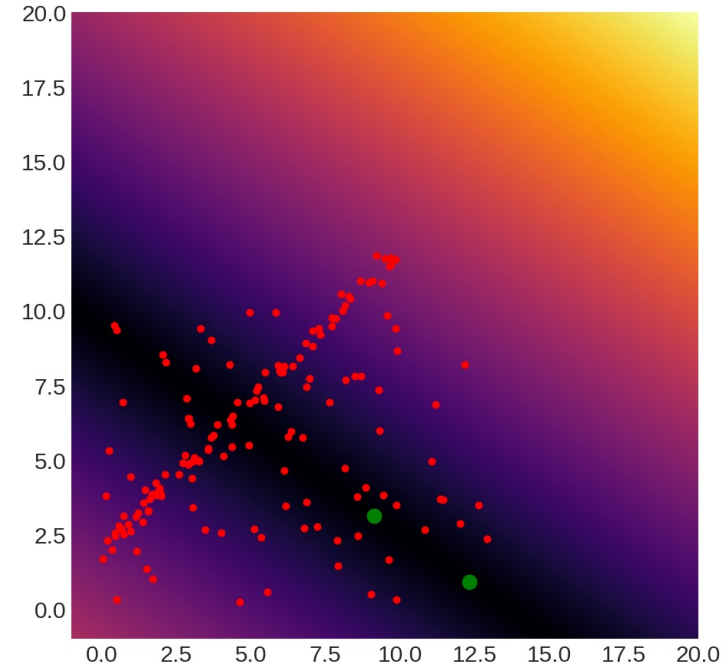    - proportional to **signed** distance to line
    - equal to signed distance if you normalize
- Generalization to N-D: hyperplane $\mathbf{w}^T\mathbf{x}+b$

# Can We Train a Network To Do It?

# Can We Train a Network To Do It?



Compute signed distance; pass through activation, e.g., ReLU

# Can We Train a Network To Do It?



Compute signed distance; pass through activation
– let's say a step function (0 or 1)

# Can We Train a Network To Do It?

# Can We Train a Network To Do It?



$\max(\mathbf{w_1}^\top \mathbf{x} + b, 0)$

$\max(\mathbf{w_1}^\top \mathbf{x} + b, 0) =$ Distance to line defined by $\mathbf{w_1}$

$x_1$

$x_2$

$\max(\mathbf{w_2}^\top \mathbf{x} + b, 0)$

$\max(\mathbf{w_2}^\top \mathbf{x} + b, 0) =$ Distance to line defined by $\mathbf{w_2}$

# Can We Train a Network To Do It?

Distance to $\mathbf{w}_1$

$x_1$

$x_2$

Next layer computes:
$$\max(\mathbf{w}_1 \text{ Distance} - \mathbf{w}_2 \text{ Distance}, 0)$$

Distance to $\mathbf{w}_2$

# Can We Train a Network To Do It?

- [Three Neurons Demo](Three Neurons Demo)

# Can We Train a Network To Do It?

*use neurons to divide your space with half-spaces (possitive value / zero)*

Result: <u>feedforward</u> neural networks with a finite number of neurons in a hidden layer can approximate any continuous function with a bounded domain
Universal Approximation Theorem.

Cybenko (1989) for neural networks with sigmoids; Hornik (1991) more generally

In practice, doesn't give a practical guarantee.
**Why?**

# Developing Intuitions

*There is no royal road to geometry. – Euclid*

- Best way: play with data, be skeptical of everything you do, be skeptical of everything you are told
- Remember: this is linear algebra, not magic
- Technique: How would you set the weights by hand if you were forced to be a deep net

# Parameters

**How many parameters does this network have?**



**Weights**: 1x2
**Parameters:** 3 (bias!)

# Parameters

**How many parameters does this network have?**



**Weights**: 1x4+4x2 = 12
**Parameters:** 12+5 = 17

# Parameters

How many parameters does this network have?



**Weights**: 3x4+4x4+4x2 = 36
**Parameters:** 36+11 = 47

# Parameters



H*P
+H

H
neurons

H*H
+H

H
neurons

H*H
+H

H
neurons

O*H
+O

O
neurons

P: 285x350 picture **(terrible!)**, H: 1000, O: 3
**102 million parameters** (400MB)

# Parameters



H neurons

- First layer converts **all** visual information into a single H dimensional vector.

- Suppose you want a neuron to represent image gradient at each pixel. **How many neurons do you need?**

- **2P** *each pixel has x and y direction gradient.*

# Parameters



Make Px1 vector

x

H*P+
H

H
neurons

H*H
+H

H
neurons

H*H
+H

H
neurons

O*H
+O

O
neurons

P: 285x350, H: 2P, O: 3
**100 billion parameters** (400GB)

# Convnets

## Keep Spatial Resolution Around

FC Neural net:
Data: vector Fx1
Transform: matrix-multiply

Convnet:
Data: image HxWxF
Transform: convolution

Make
Px1
vector

x

the spacial information
is ignored.

Keep
Image
Dims

# Convnet (Feature Volume)



Height: 300
Width: 500
Depth: 3

Height: 32
Width: 32
Depth: 3

Height

Width

Depth
or
Channel

# Convnet

**Fully connected:**
Connects to everything

32
neuron
32
3

**Convnet:**
Connects locally
Filter Small & Reused

32
neuron
32
3

# Convnet

## Neuron is the same: weighted linear average



32

neuron

32

3

$F_w$

$F_h$

c

$$\sum_{i=1}^{F_h} \sum_{j=1}^{F_w} \sum_{k=1}^{c} F_{i,j,k} I_{y+i,x+j,k}$$

# Convnet

## Neuron is the same: weighted linear average

$F_w$

3_

32

3

Filter is local in space: sum only over $F_h \times F_w$ pixels

Filter is global over channels/depth: sum over all channels

$$\sum_{i=1}^{F_h} \sum_{j=1}^{F_w} \sum_{k=1}^{c} F_{i,j,k} I_{y+i,x+j,k}$$

# Convnet

## Get spatial output by sliding filter over image



$$\sum_{i=1}^{F_h}\sum_{j=1}^{F_w}\sum_{k=1}^{c} F_{i,j,k} I_{y+i,x+j,k}$$

# Differences From Earlier Filtering

(a) #input channels can be greater than one
(b) forget you learned the difference between convolution and cross-correlation

| I11 | F11 | F12 | F13 | I15 | I16 |
|-----|-----|-----|-----|-----|-----|
| I21 | F21 | F22 | F23 | I25 | I26 |
| I31 | F31 | F32 | F33 | I35 | I36 |
| I41 | I42 | I43 | I44 | I45 | I46 |
| I51 | I52 | I53 | I54 | I55 | I56 |

Output[1,2]
= I[1,2]*F[1,1] + I[1,3]*F[1,2]
+ … + I[3,4]*F[3,3]

# Convnet

**How big is the output?**

32

5

5

32

3

Height? 32-5+1=28
Width? 32-5+1=28
Channels? 1

One filter not very useful by itself

# Multiple Filters

You've already seen this before

Input:
400x600x1

Output:
400x600x2

# Convnet
## Multiple out channels via multiple filters.
## **How big is the output?**



Depth
Dimension

32

5

5

200

200 Filters

32

3

Height? 32-5+1=28
Width? 32-5+1=28
Channels? 200

# Convnet
## Multiple out channels via multiple filters.
## How big is the output?

5

32

5

5

32

3

Height? 32-5+1=28
Width? 32-5+1=28
Channels? 200

# Convnet, Summarized

Neural net:
series of matrix-multiplies
parameterized by **W**,**b** +
nonlinearity/activation
Fit by gradient descent

Convnet:
series of convolutions
parameterized by **F,b** +
nonlinearity/activation
Fit by gradient descent

# One Additional Subtlety – Stride

## Warmup: how big is the output spatially?

| | | | | | | |
|---|---|---|---|---|---|---|
| F11 | F12 | F13 | I14 | I15 | I16 | I17 |
| F21 | F22 | F23 | I24 | I25 | I26 | I27 |
| F31 | F32 | F33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Normal (Stride 1):
5x5 output

Example credit: Karpathy and Fei-Fei

# One Additional Subtlety – Stride

## Stride: skip a few (here 2)

| | | | | | | |
|---|---|---|---|---|---|---|
| F11 | F12 | F13 | I14 | I15 | I16 | I17 |
| F21 | F22 | F23 | I24 | I25 | I26 | I27 |
| F31 | F32 | F33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

## Normal (Stride 1): 5x5 output

# One Additional Subtlety – Stride

## Stride: skip a few (here 2)

| | | | | | | |
|------|------|------|------|------|------|------|
| I11 | I12 | F11 | F12 | F13 | I16 | I17 |
| I21 | I22 | F21 | F22 | F23 | I26 | I27 |
| I31 | I32 | F31 | F32 | F33 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

## Normal (Stride 1):
## 5x5 output

Example credit: Karpathy and Fei-Fei

# One Additional Subtlety – Stride

## Stride: skip a few (here 2)

| | | | | | | |
|---|---|---|---|---|---|---|
| I11 | I12 | I13 | I14 | F11 | F12 | F13 |
| I21 | I22 | I23 | I24 | F21 | F22 | F23 |
| I31 | I32 | I33 | I34 | F31 | F32 | F33 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

### Normal (Stride 1): 5x5 output

### Stride 2 convolution: 3x3 output

# One Additional Subtlety – Stride

## What about stride 3?

| | | | | | | |
|---|---|---|---|---|---|---|
| F11 | F12 | F13 | I14 | I15 | I16 | I17 |
| F21 | F22 | F23 | I24 | I25 | I26 | I27 |
| F31 | F32 | F33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Normal (Stride 1):
5x5 output

Stride 2 convolution:
3x3 output

Example credit: Karpathy and Fei-Fei

# One Additional Subtlety – Stride

## What about stride 3?

| | | | | | | |
|---|---|---|---|---|---|---|
| I11 | I12 | I13 | F11 | F12 | F13 | I17 |
| I21 | I22 | I23 | F21 | F22 | F23 | I27 |
| I31 | I32 | I33 | F31 | F32 | F33 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Normal (Stride 1):
5x5 output

Stride 2 convolution:
3x3 output

Stride 3 convolution:
Not Common!

Example credit: Karpathy and Fei-Fei

# One Additional Subtlety

Zero padding is extremely common, although other forms of padding do happen

$g$ | $? ? ? ?$ | $g$

$f$

$g$ | | $g$

Symm: fold sides over

Circular/Wrap: wrap around

pad/fill: add value, often 0

# In General



Output Size

$$\left\lfloor \frac{(N - F)}{S} \right\rfloor + 1$$

# In General



N+2P

F

F

S

N+2P

Output Size
with Padding

$$\left\lfloor \frac{(N + 2P - F)}{S} \right\rfloor + 1$$

# In General



N+2P

N+2P

F

F

S

With Proper Padding:
Output Size Preserved

$$\left\lfloor \frac{N}{S} \right\rfloor$$

Proper Padding
Size?

$$\left\lfloor \frac{F}{2} \right\rfloor$$

# More Examples

Input volume: **32x32x3**
Receptive fields: **5x5**, **stride 1**
Number of neurons: **5**

$$\frac{(N - F)}{s} + 1$$



# Output volume size?

# More Examples

Input volume: **32x32x3**
Receptive fields: **5x5**, **stride 1**
Number of neurons: **5**

$$\frac{(N - F)}{s} + 1$$

Output volume: (32 - 5) / 1 + 1 = 28, so: **28x28x5**

# Number of Parameters?

bias

5x5x3 × 5 +5
each filter   we have five filters

# More Examples

Input volume: **32x32x3**
Receptive fields: **5x5**, **stride 1**
Number of neurons: **5**

$$\frac{(N - F)}{s} + 1$$



Output volume: (32 - 5) / 1 + 1 = 28, so: **28x28x5**
How many parameters? **5x5x3x5 + 5 = 380**

Smaller than FC (32x32x3x5+5)=15365

# More Examples

Input volume: **32x32x3**
Receptive fields: **5x5**, **stride** <span style="color:red">**3**</span>
Number of neurons: **5**

$$\frac{(N - F)}{s} + 1$$



# Output volume size?

# More Examples

Input volume: **32x32x3**
Receptive fields: **5x5**, **stride 3**
Number of neurons: **5**



Output volume: (32 - 5) / 3 + 1 = 10, so: **10x10x5**

# Number of Parameters?

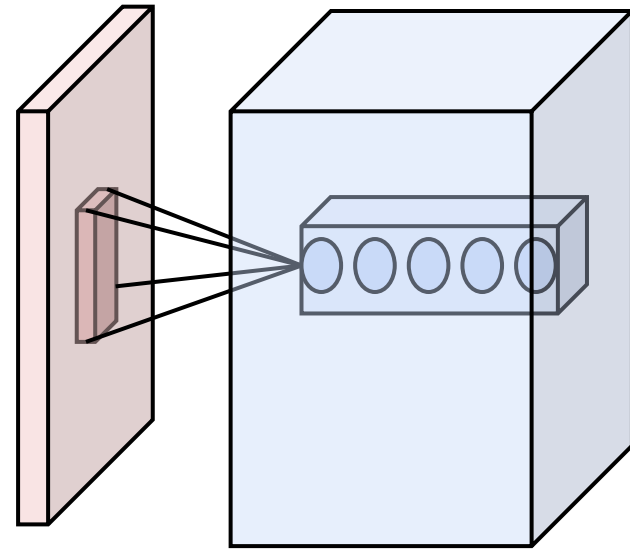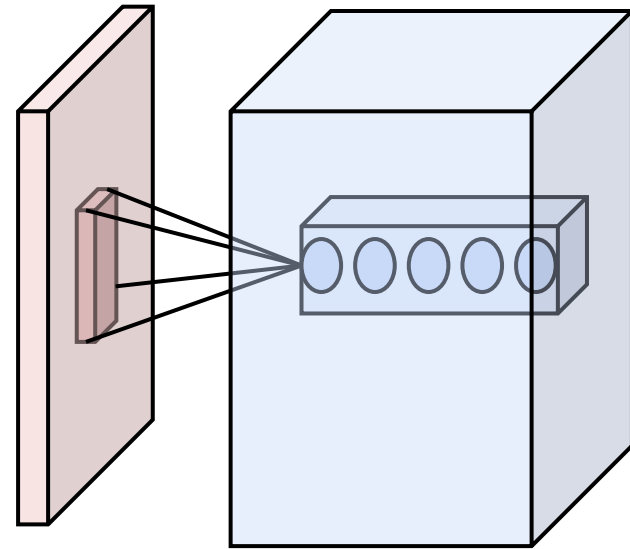# More Examples

Input volume: **32x32x3**
Receptive fields: **5x5**, **stride 3**
Number of neurons: **5**



Output volume: (32 - 5) / 3 + 1 = 10, so: **10x10x5**
How many parameters? **5x5x3x5 + 5 = 380. Same!**

# Other Layers – Pooling

Idea: want spatial resolution of activations / images smaller; applied per-channel



Max-pool
2x2 Filter
Stride 2

# Other Layers – Pooling

Idea: want spatial resolution of activations / images smaller; applied per-channel

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 1 | 3 | 4 |

Avg-pool
2x2 Filter
Stride 2

→

| 3.25 | 5.25 |
|------|------|
| 1.75 | 2.0  |

# Other Layers – Pooling

Idea: want spatial resolution of activations / images smaller; applied per-channel

| I11 | I12 | I13 | I14 | I15 | I16 | I17 |
|-----|-----|-----|-----|-----|-----|-----|
| I21 | I22 | I23 | I24 | I25 | I26 | I27 |
| I31 | I32 | I33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Max-pool
3x3 Filter
Stride 2

O11

O11 = maximum value in blue box

# Other Layers – Pooling

Idea: want spatial resolution of activations / images smaller; applied per-channel

| | | | | | | |
|---|---|---|---|---|---|---|
| I11 | I12 | I13 | I14 | I15 | I16 | I17 |
| I21 | I22 | I23 | I24 | I25 | I26 | I27 |
| I31 | I32 | I33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Max-pool
3x3 Filter
Stride 2

| | |
|---|---|
| O11 | O12 |

O12 = maximum value in blue box

# Other Layers – Pooling

Idea: want spatial resolution of activations / images smaller; applied per-channel

| I11 | I12 | I13 | I14 | I15 | I16 | I17 |
|-----|-----|-----|-----|-----|-----|-----|
| I21 | I22 | I23 | I24 | I25 | I26 | I27 |
| I31 | I32 | I33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Max-pool
3x3 Filter
Stride 2

| O11 | O12 | O13 |
|-----|-----|-----|

O13 = maximum value in blue box

# Other Layers – Pooling

Idea: want spatial resolution of activations / images smaller; applied per-channel

| | | | | | | |
|---|---|---|---|---|---|---|
| I11 | I12 | I13 | I14 | I15 | I16 | I17 |
| I21 | I22 | I23 | I24 | I25 | I26 | I27 |
| I31 | I32 | I33 | I34 | I35 | I36 | I37 |
| I41 | I42 | I43 | I44 | I45 | I46 | I47 |
| I51 | I52 | I53 | I54 | I55 | I56 | I57 |
| I61 | I62 | I63 | I64 | I65 | I66 | I67 |
| I71 | I72 | I73 | I74 | I75 | I76 | I77 |

Max-pool
3x3 Filter
Stride 2

$\longrightarrow$

| | | |
|---|---|---|
| O11 | O12 | O13 |
| O21 | O22 | O23 |
| O31 | O32 | O33 |

Smaller Resolution→Smaller Computation

# Squeezing a Loaf of Bread



Max-pool
2x2 Filter
Stride 2

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 1 | 3 | 4 |

| 6 | 8 |
|---|---|
| 3 | 4 |

# Example Network

Suppose we want to convert a 32x32x3 image into a 10x1 vector of classification results



Figure Credit: Karpathy and Fei-Fei; see http://cs231n.stanford.edu/

# Example Network

input: [32x32x3]
CONV with 10 3x3 filters, stride 1, pad 1:
gives: [32x32x10]
new parameters: (3*3*3)*10 + 10 = 280
RELU
CONV with 10 3x3 filters, stride 1, pad 1:
gives: [32x32x10]
new parameters: (3*3*10)*10 + 10 = 910
RELU
POOL with 2x2 filters, stride 2:
gives: [16x16x10]
parameters: 0

Slide credit: Karpathy and Fei-Fei

# Example Network

Previous output: [16x16x10]
CONV with 10 3x3 filters, stride 1:
gives: [16x16x10]
new parameters: (3*3*10)*10 + 10 = 910
RELU
CONV with 10 3x3 filters, stride 1:
gives: [16x16x10]
new parameters: (3*3*10)*10 + 10 = 910
RELU
POOL with 2x2 filters, stride 2:
gives: [8x8x10]
parameters: 0

# Example Network

Conv, Relu, Conv, Relu, Pool continues until it's [4x4x10]

Fully-Connected FC layer to 10 neurons
(which are our class scores)
Number of parameters:
10 * 4 * 4 * 10 + 10 = 1610

done!

# An Alternate Conclusion

Conv, Relu, Conv, Relu, Pool continues until it's [4x4x10]
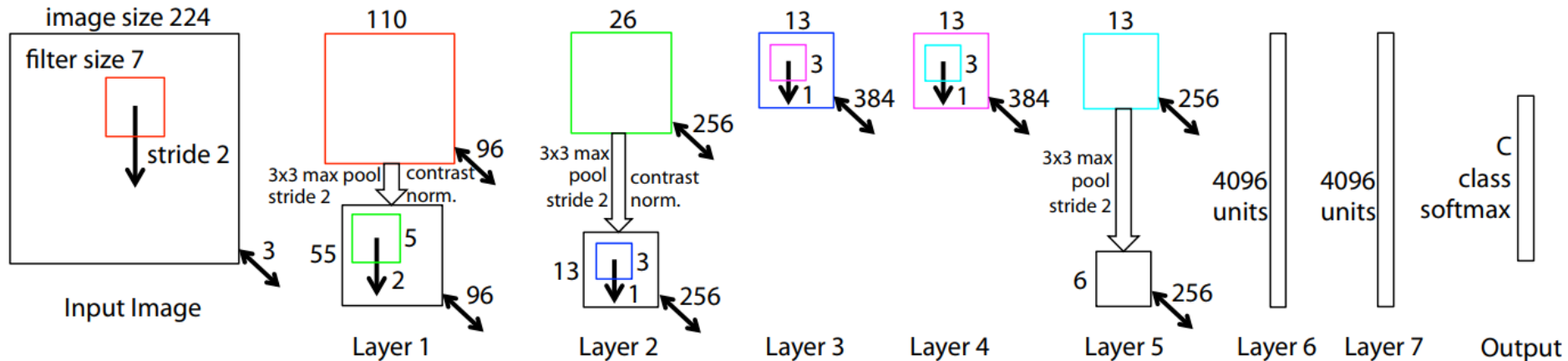
Average POOL 4x4x10 to 10 neurons
Fully-Connected FC layer to 10 neurons
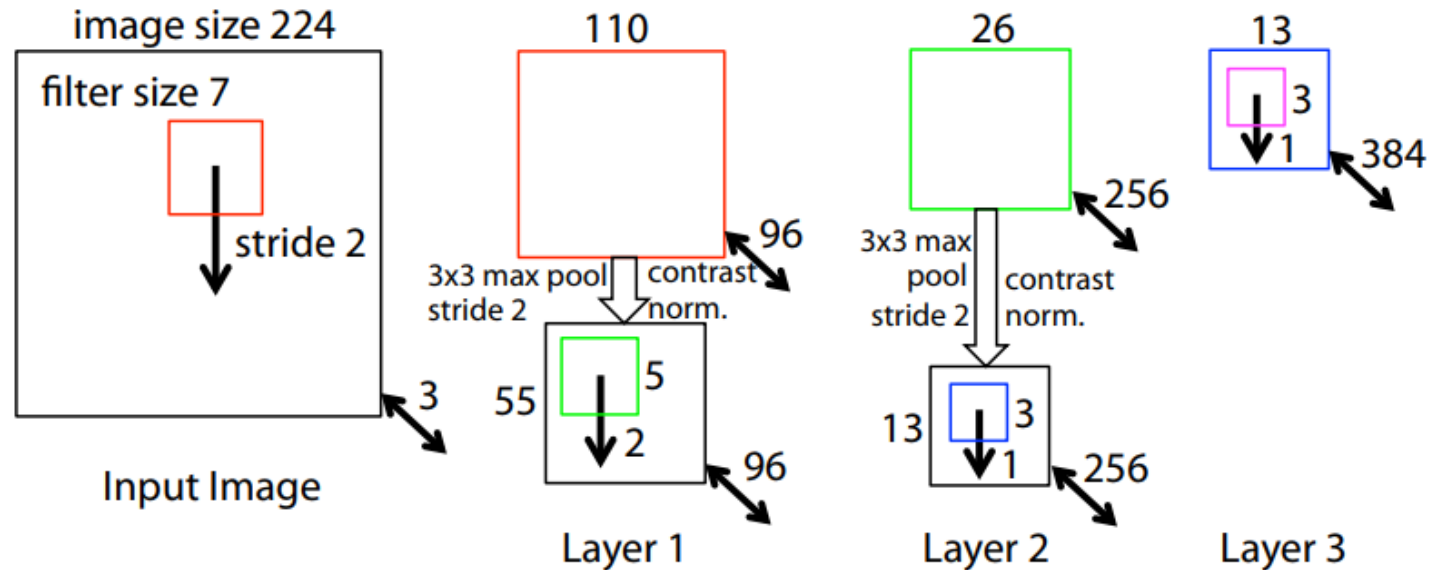(which are our class scores)
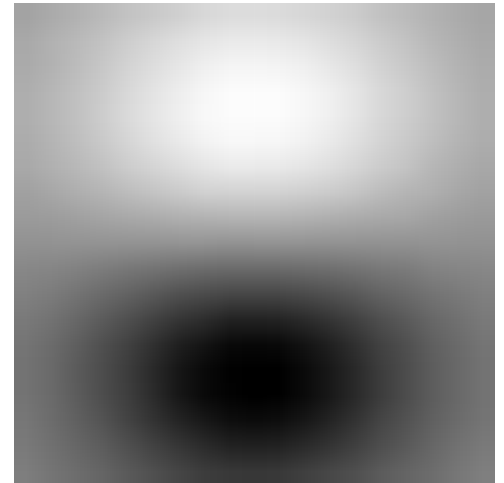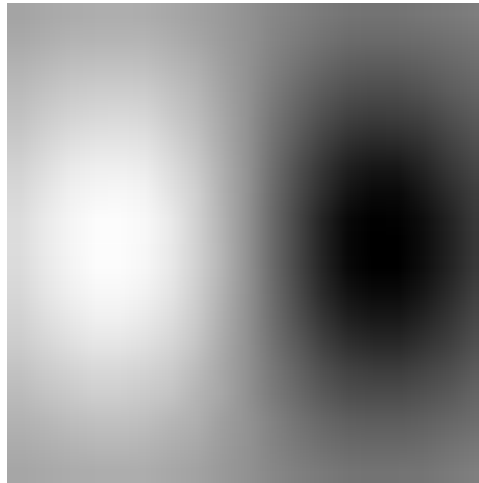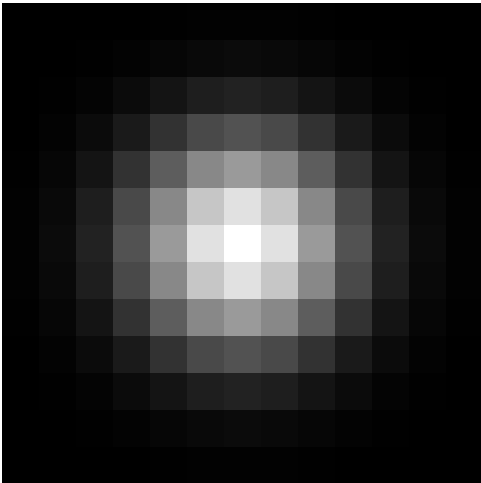Number of parameters:
10 * 10 + 10 = 110

done!

# Example Network



Figure Credit: Zeiler and Fergus, Visualizing and Understanding Convolutional Networks. ECCV 2014

# Example Network



image size 224
filter size 7
stride 2
Input Image

110
3x3 max pool stride 2
contrast norm.
96

55
5
2
96
Layer 1

26
3x3 max pool stride 2
contrast norm.
256

13
3
1
256
Layer 2

13
3
1
384
Layer 3

3

(1) filter image with 96 7x7 filters with stride 2
(2) ReLU
(3) 3x3 max pool with stride 2 (and *contrast normalization – now ignored)*

Figure Credit: Zeiler and Fergus, Visualizing and Understanding Convolutional Networks. ECCV 2014

# What Do The Filters Represent?

Recall: filters are images and we can look at them

# What Do The Filters Represent?



First layer filters (11x11) of a network (AlexNet) trained to distinguish 1000 categories of objects (ImageNet)

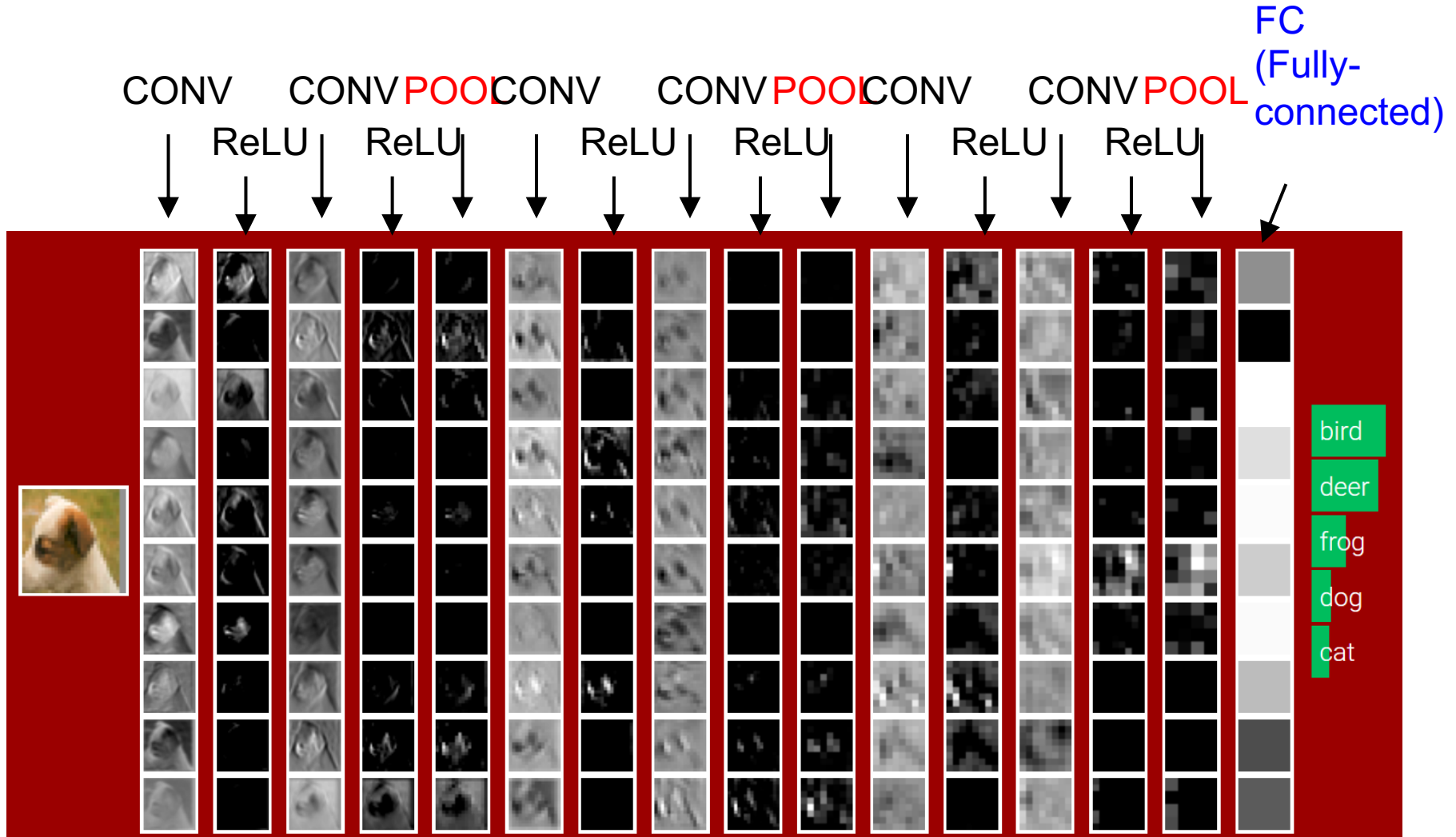Remember these filters go over color.

For the interested: Gabor filter

Figure Credit: Karpathy and Fei-Fei

# What Do The Filters Do?



Figure Credit: Karpathy and Fei-Fei; see http://cs231n.stanford.edu/

# Next Class: More CNNs